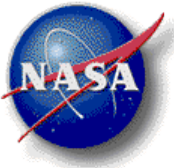
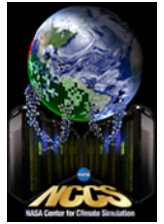


NCCS User Forum

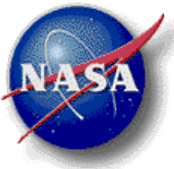
March 20, 2012



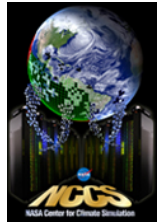
Agenda – March 20, 2012



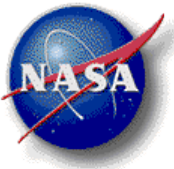
- Welcome & Introduction (Lynn Parnell)
- Discover Update (Dan Duffy)
- Dirac DMF Growth Rate – Issue & Remediation (Tom Schardt)
- NCCS Operations & User Services (Ellen Salmon)
- Question & Answer
- Breakout Sessions
 - Overview of Debugging on Discover (Doris Pan)
 - OpenCL Framework for Heterogeneous CPU/GPU Programming (George Fekete)



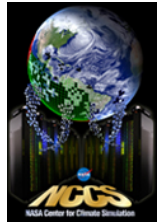
Welcome & Intro - Announcements



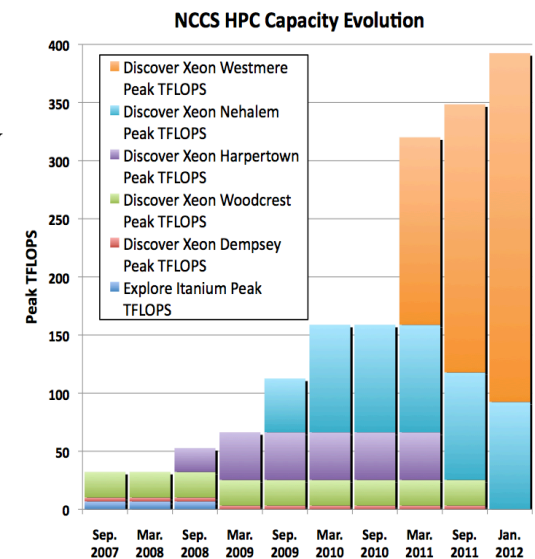
- Two Breakouts after main session's Q&A:
 - “Overview of Debugging on Discover” (Doris Pan)
 - “OpenCL Framework for Heterogeneous CPU/GPU Programming” (George Fekete)
- Sign-in sheet
 - Column for concerns & suggestions for topics for Brown Bag seminars or future NCCS User Forums
- SMD Allocation Applications Deadline: 20 March (today!)
 - <http://www.hec.nasa.gov/request/announcements.html#spring2012>
- New employees: Doris Pan, Bill Woodford, Ben Bledsoe

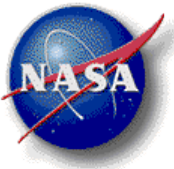


Accomplishments

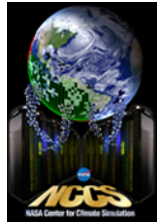


- Discover Westmere replacement for old cores (SCUs 1-4)
 - Adds 50 TFLOPs capacity
- Dali Analysis additions (9 systems), each with:
 - Westmere nodes (dual-socket, hex-core), 192 GB memory
 - 2 NVIDIA “Tesla” general-purpose Graphical Processing Units
- Archive upgrades
 - New tape library & higher capacity media/tape drives
 - Additional I/O servers for higher bandwidth
 - High availability/failover for NFS exports to Discover
- Data Services
 - Controlled data-sharing with designated colleagues (e.g., Land Information Systems)
 - Ramp-up of Earth System Grid supporting IPCC AR5 (40 TB downloaded so far)
- NCCS staff participated in SMD’s testing of Goddard’s Nebula cloud services system.

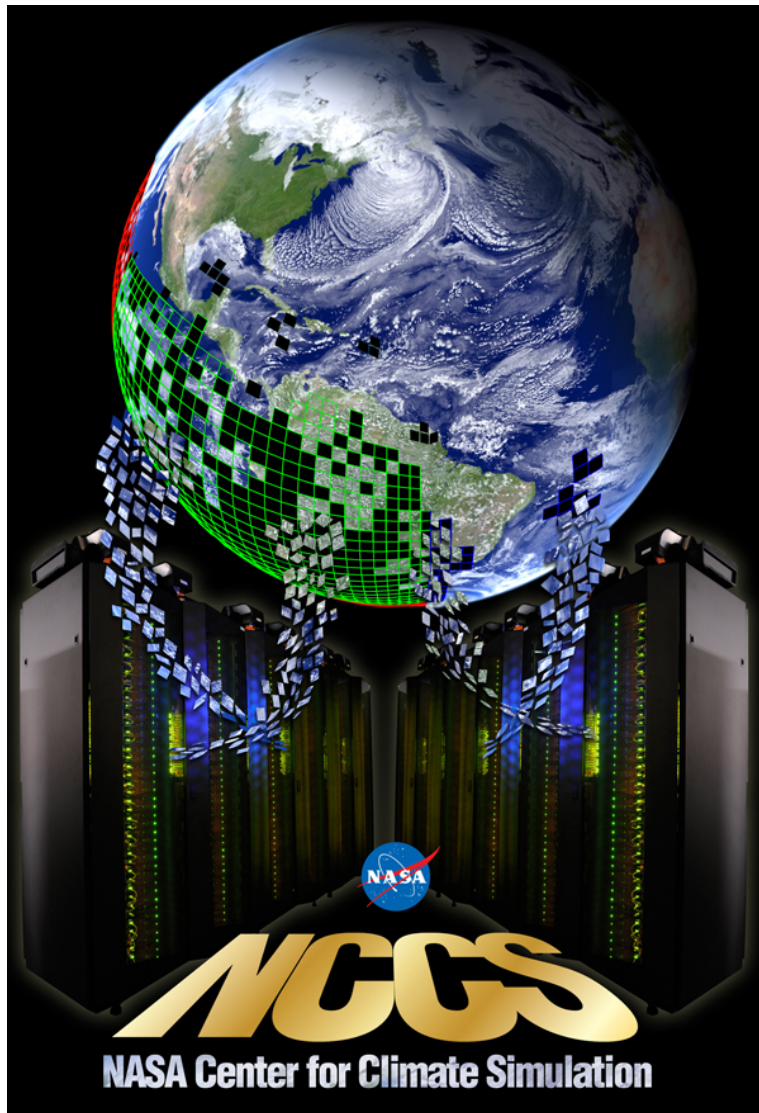
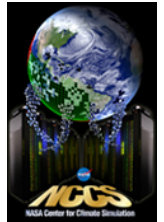




Agenda – March 20, 2012



- Welcome & Introduction (Lynn Parnell)
- Discover Updates (Dan Duffy)
- Dirac DMF Growth Rate – Issue & Remediation (Tom Schardt)
- NCCS Operations & User Services (Ellen Salmon)
- Question & Answer
- Breakout Sessions
 - Overview of Debugging on Discover (Doris Pan)
 - OpenCL Framework for Heterogeneous CPU/GPU Programming (George Fekete)



Discover Updates

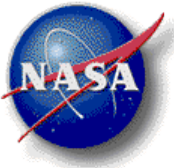
Chilled Water Outage

Recent Upgrades

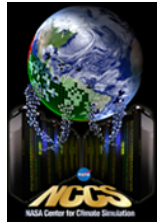
SCU8 Plans (Diskless?)

Sandy Bridge Performance

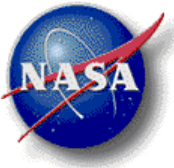
March 20, 2012



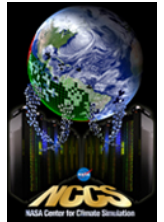
February Chilled Water Outage



- NCCS continuity of operations (no downtime) was considered highly desirable
- Attempt to continue minimal operations during chilled water outage failed
 - All NCCS services went down for the day
 - Multiple disk drives in Discover cluster, archive failed immediately
 - Continued disk failures plague the cluster & archive even now. Total ~50
 - Recovery work continues
- Disk failures indicate systemic reliability problem in some NCCS disks
 - Temperature rise for chilled water outage didn't appear to exceed operation specs
 - Vendors & manufacturers are being engaged to solve reliability problem
 - Specific, limited device problems have been noted
 - Root cause analysis & update of procedures are planned
 - Details of analysis will be provided when complete.



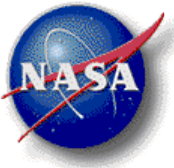
Discover and Dali Upgrades Since July 2011



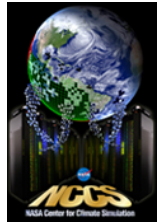
- SCUs 1-4 Xeon Westmere upgrades, August 2011 through January 2012
- Swapped 1,032 old Dell & IBM nodes (circa 2006-2008) for 1,032 new IBM Westmere nodes
- Re-used racks & switches
- Net addition of ~3M SBU per year

Discover Feature	Old SCU1&2	Old SCU3&4	New SCUs 1-4
Processor Type	Woodcrest	Harpertown	Westmere
Cores per node	4	8	12
SBU rate per node	0.20	0.38	0.95

- Nine new Dali Analysis Westmere systems, each with:
 - Dual socket, hex-core Intel Westmere 2.8 GHz cores
 - 192 GB memory; QDR Infiniband
 - Two NVIDIA M2070 Tesla GPUs
 - 10 GbE Network Interface Cards
 - All GPFS nobackup file systems mounted
- Using “proxy” Bastion service, accessible by direct login:
 - ssh dali-gpu.nccs.nasa.gov (new nodes)
 - ssh dali.nccs.nasa.gov – might land on a new node with a GPU (round robin)
 - Details at <http://www.nccs.nasa.gov/primer/getstarted.html#login>

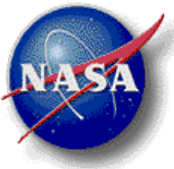


Upgrades Later This Year

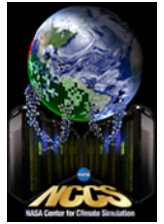


- Compute Upgrade
 - NCCS is designing and procuring the next compute upgrade
 - Similar to the current architecture
 - Intel Xeon Sandy Bridge cores (16 per node), InfiniBand, GPFS, PBS
 - Strong preference for a “containerized” solution
 - Thinking about using “diskless” nodes (*next slide*)
 - No node-local disk, but GPFS file systems are mounted
- Discover Storage Upgrades
 - GPFS storage capacity increase ~1 PB
 - Prototype file system for many small files





SCU8: Diskless Nodes Discussion



Challenges

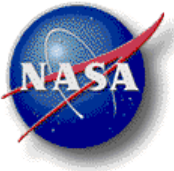
- No swap space
 - Might need code changes
 - *Would nodes need more physical memory?*
- No local scratch space
- May need to manage two sets of codes, for nodes with and without local disk

Advantages

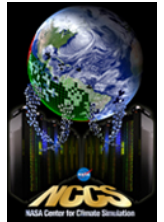
- Uses less power and cooling
- Less expensive
- Pleaides at NAS is diskless, so many applications are already ported & needed code changes are well known
- Majority of (current) Discover nodes will still have local disk (can specify to PBS)

Question for you:

What impact would diskless nodes have on your research?



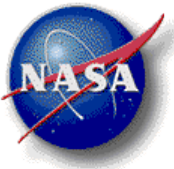
Intel Sandy Bridge Processors Performance



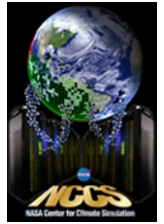
Benchmark	Westmere	Sandy Bridge	Speedup
ModelE (1)	722.33	494.21	1.46
ModelE (2)	412.14	310.33	1.33
Cubed Sphere (3)	724.13	425.68	1.72
WRF (4)	43.45	51.95	1.63

Notes:

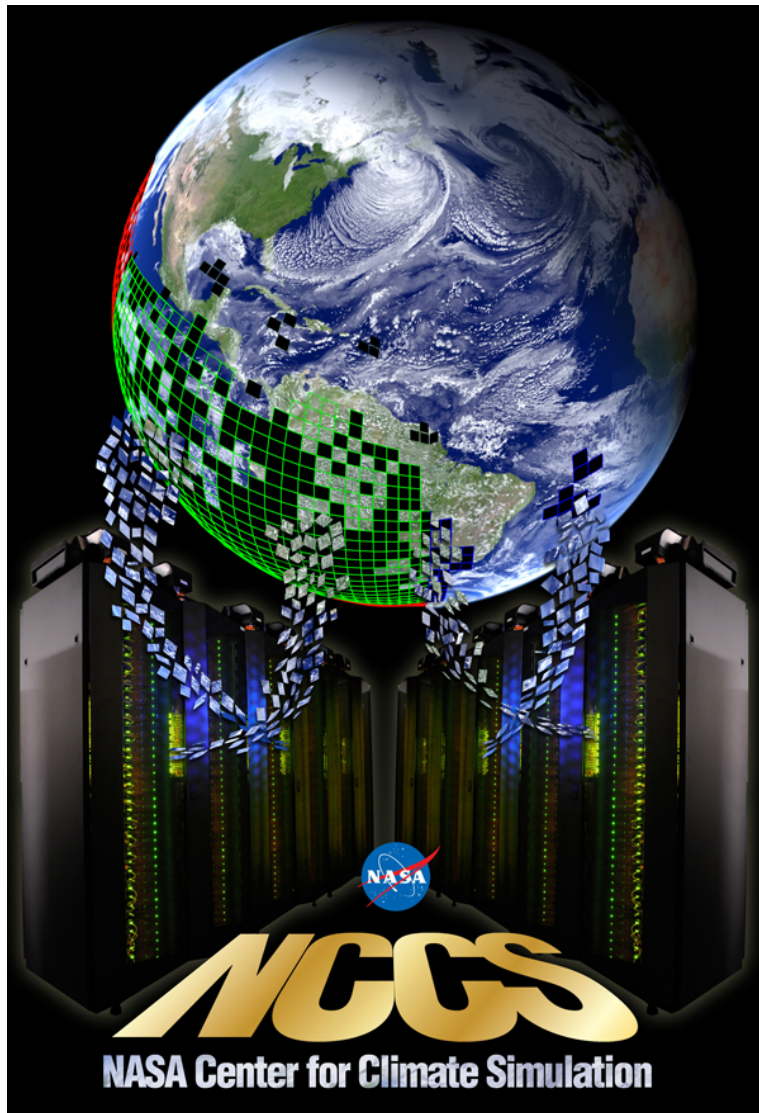
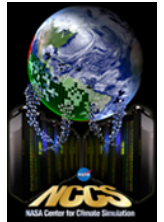
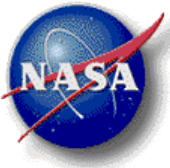
- (1) E4TcadF40 rundeck; 8 MPI processes; $[\text{time}(5\text{-day}) - \text{time}(1\text{day})] / 4$
- (2) E4TcadF40 rundeck; 16 MPI processes; $[\text{time}(5\text{-day}) - \text{time}(1\text{day})] / 4$; saturated SandyBridge node
- (3) Cubed Sphere; benchmark 1; 12 MPI processes; saturated discover node
- (4) Snow storm of Jan. 23-24, 2005; 15 km horizontal grid (56x92), 38 vertical levels, 24 hour run



Agenda – March 20, 2012

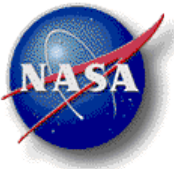


- Welcome & Introduction (Lynn Parnell)
- Discover Updates (Dan Duffy)
- Dirac DMF Growth Rate – Issue & Remediation (Tom Schardt)
- NCCS Operations & User Services (Ellen Salmon)
- Question & Answer
- Breakout Sessions
 - Overview of Debugging on Discover (Doris Pan)
 - OpenCL Framework for Heterogeneous CPU/GPU Programming (George Fekete)

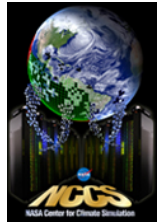


Dirac DMF Growth Rate – Issue & Remediation

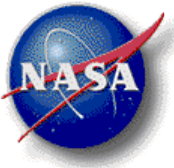
March 20, 2012



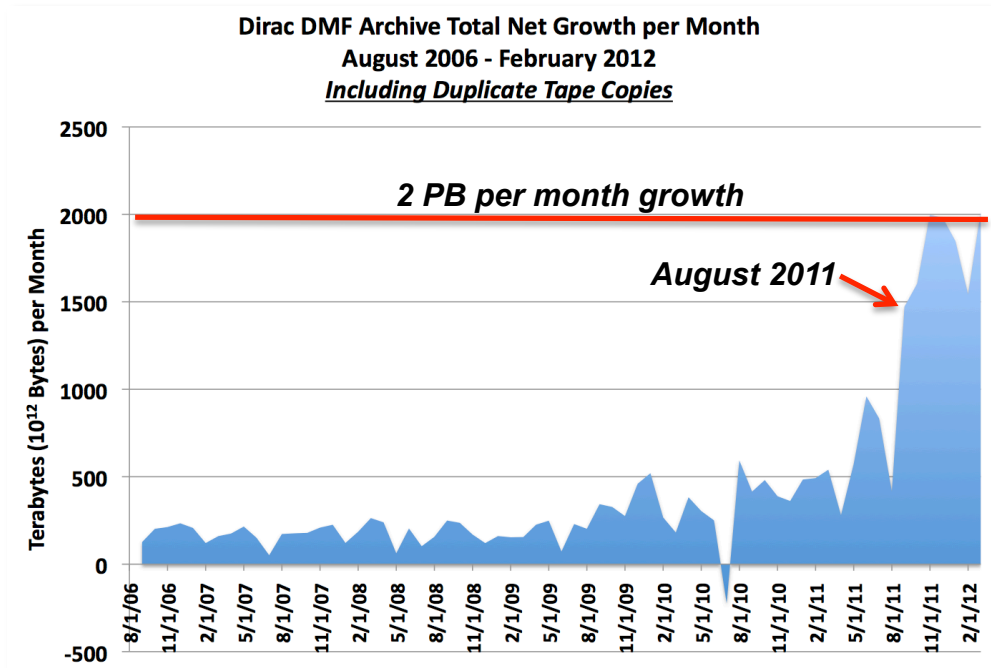
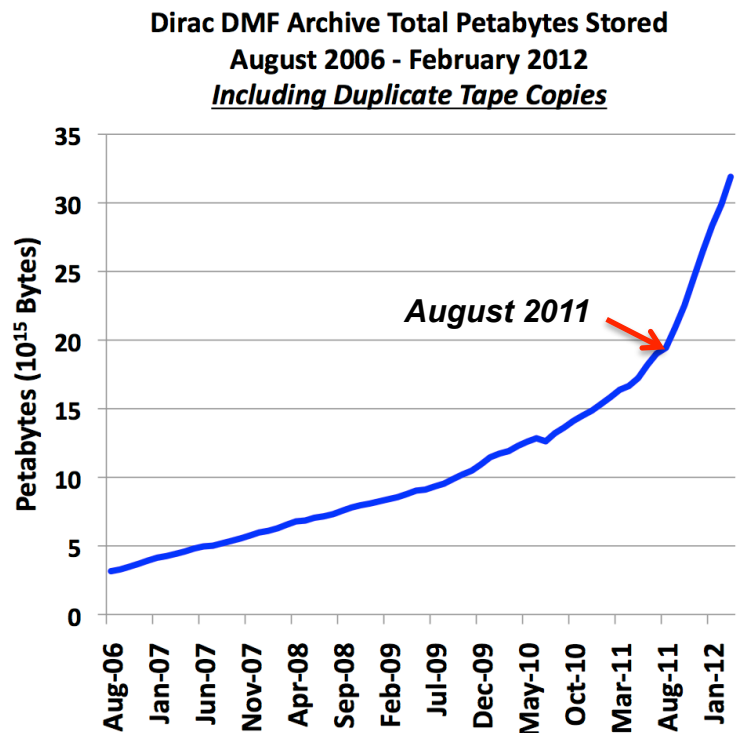
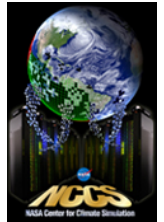
Current Dirac DMF Growth Rate Is Unsupportable with Two Tape Copies



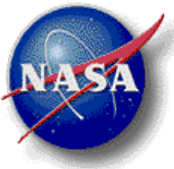
- Over the last 7 months, we have had a dramatic increase in archive growth.
- FY12 cost to retain 2 tape copies at current growth rate:
 - \$2.3M (new data: \$1.8M; migration from oldest tapes: \$0.5M)
 - Current growth rate: 2 Petabytes/month including both tape copies.
- By reducing archive to one tape copy, NCCS could free an additional \$1M to continue to expand the Discover cluster to better satisfy the ever-growing demand for processing.
- Thus NCCS is strongly considering making only one tape copy of new archive data as of May 1, 2012.
- ***Specific high-value data could still get two archive tape copies.***
- If the single tape copy approach is adopted, over time the second tape copies of most pre-existing data would also be removed after verifying the tape containing the remaining copy is readable.



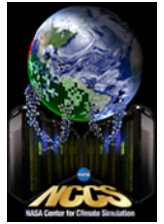
Dirac DMF Archive Total Stored & Growth



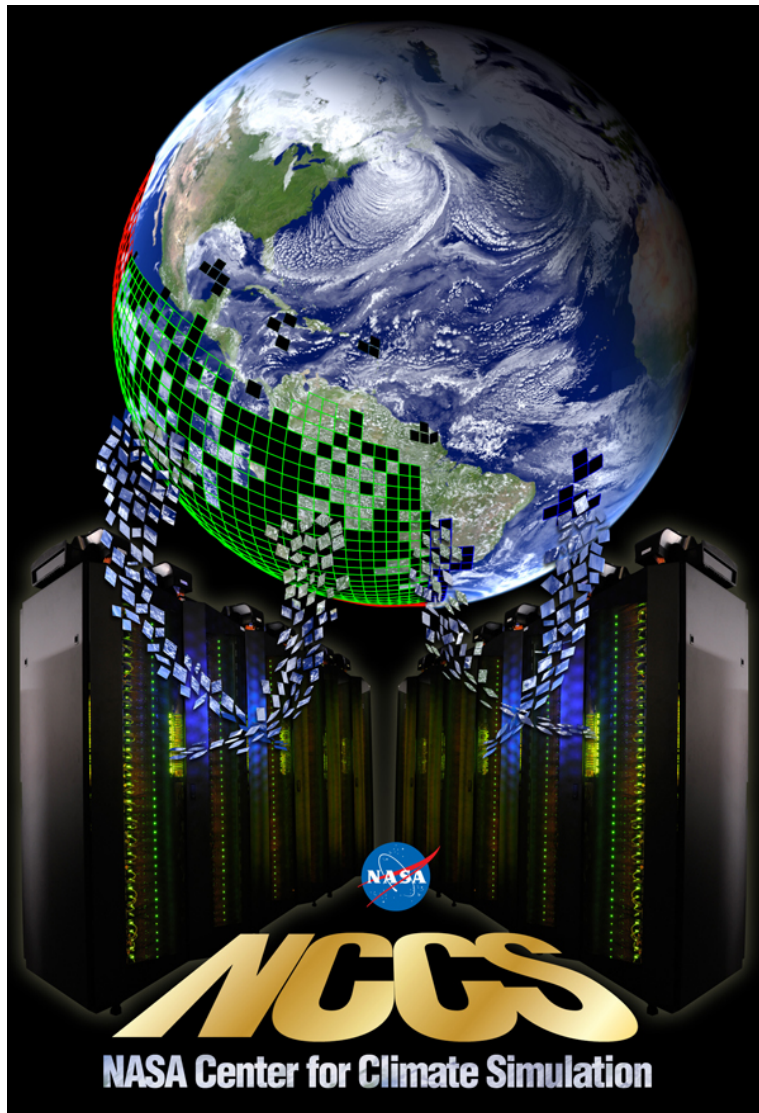
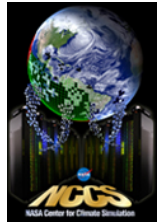
- August 2011: dramatic increase in archive growth.
 - Charts above include duplicate tape copies



Agenda – March 20, 2012



- Welcome & Introduction (Lynn Parnell)
- Discover Update (Dan Duffy)
- Dirac DMF Growth Rate – Issue & Remediation (Tom Schardt)
- NCCS Operations & User Services (Ellen Salmon)
- Question & Answer
- Breakout Sessions
 - Overview of Debugging on Discover (Doris Pan)
 - OpenCL Framework for Heterogeneous CPU/GPU Programming (George Fekete)



HPC Operations and User Services

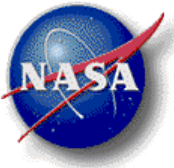
PBS 11 & SLES 11 SP 1

Login Bastion Services

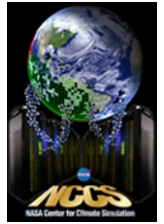
NCCS Primer

NCCS Brown Bag Seminars

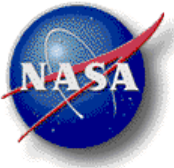
March 20, 2012



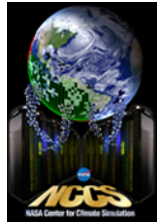
Coming Discover Changes: Discover PBS 11 and Linux SLES 11 SP1



- Reason for upgrades:
 - PBS 11 will get the bugfix for excruciatingly long PBS startups
 - SLES 11 SP1 is required to maintain current Linux security patches
- First will be tested extensively by NCCS staff, then pioneer users
 - PBS 11: no script changes should be needed, per vendor Altair
 - SLES 11 SP 1: we'll test for, and document, needed changes due to updated libraries and Linux kernel
- Planning for phased, rolling deployments with minimal downtime
- PBS 11 efforts will begin in April
- SLES 11 SP1 will follow later



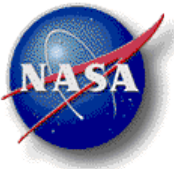
NCCS Bastion Service



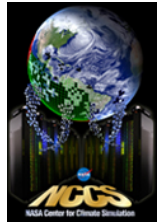
- Arbitrates SSH login access to NCCS hosts.
- Redundant service supported by multiple physical systems.
- Two modes: “Bastion” and “Proxy” (new).
- Already used for Discover logins.
- **Becomes the default mechanism for SSH logins to Dirac and Dali on March 26, 2012.**
- Why the change? Improves protection and access management, with little-to-no change for users.

Q: Can I still transfer files directly into Dirac and Dali?

A: Yes, following a few one-time-setup steps.



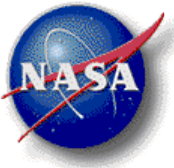
NCCS Bastion Service (continued)



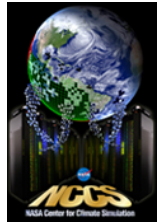
- Bastion Mode:
 - Best suited for SSH remote shell, i.e.: command line access.
 - Used exactly the same way as existing login.nccs.nasa.gov service.
- Proxy Mode (new):
 - • Best suited for SSH **file copy (scp or sftp)** access to NCCS hosts.
 - Submit ticket to request the required LDAP user account attribute.
 - Also requires one-time client-side configuration:
 - A few lines in your .ssh/config file.

Please see the NCCS Primer, System Login, for more details:

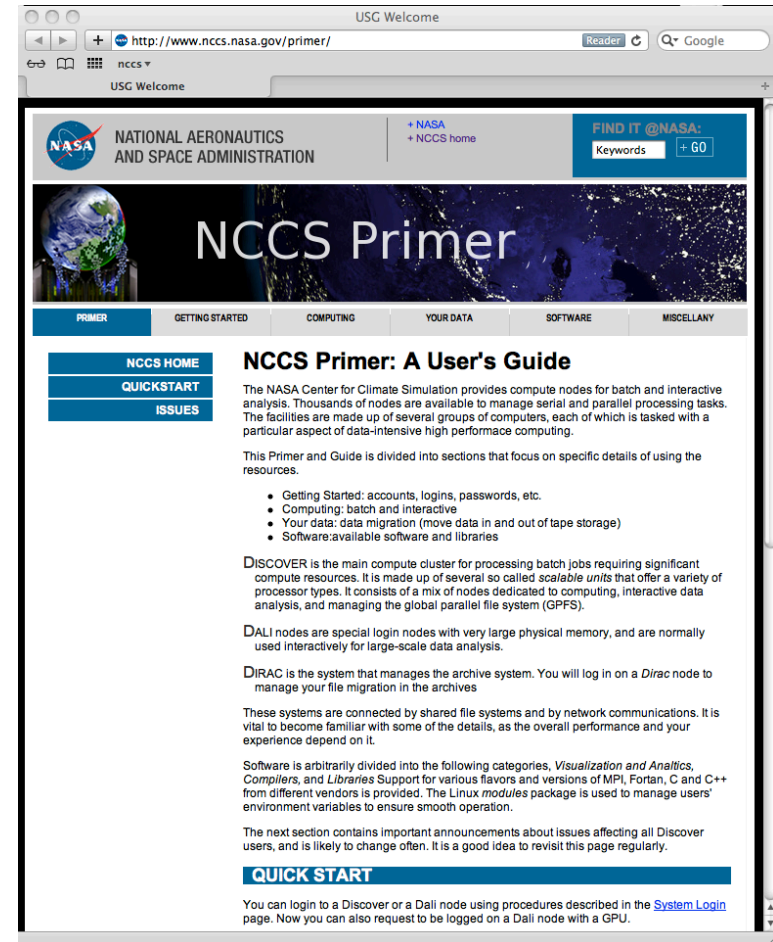
<http://www.nccs.nasa.gov/primer/getstarted.html#login>



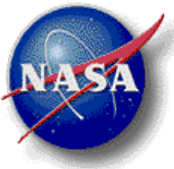
New NCCS Primer Web Pages



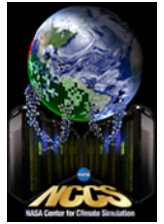
- Significant update to NCCS web pages
- Updates to NCCS basic content
- Detailed info from extensive ticket review
- Tabs at top to navigate to major categories
- Please suggest additions via email to support@nccs.nasa.gov



Initial page of the new NCCS Primer.



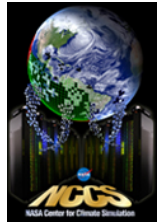
New: Twice-Monthly NCCS Brown Bag Seminars



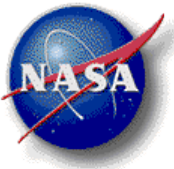
- In GSFC Building 33 (as available)
- Content will be available on the NCCS web site
- Suggest topics of interest on today's signup sheet, or via email to support@nccs.nasa.gov
- What day(s) of the week will work best for you?



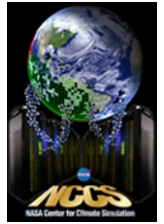
Twice-Monthly NCCS Brown Bag Seminars: Some Proposed Topics



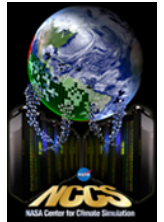
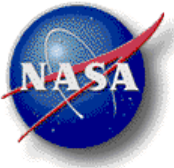
- NCCS Storage Usage & Best Practices Tips
- Using Totalview on Discover
- Monitoring PBS Jobs and Memory
- Introduction to Using Matlab with GPUs
- Best Practices for Using Matlab with GPUs
- Introduction to the NCCS Discover Environment
- Scientific Computing with Python
- Using CUDA with NVIDIA GPUs
- Using OpenCL
- Using GNU Octave
- Using Database Filesystems for Many Small Files



Agenda – March 20, 2012



- Welcome & Introduction (Lynn Parnell)
- Discover Update (Dan Duffy)
- Dirac DMF Growth Rate – Issue & Remediation (Tom Schardt)
- NCCS Operations & User Services (Ellen Salmon)
- Question & Answer
- Breakout Sessions
 - Overview of Debugging on Discover (Doris Pan)
 - OpenCL Framework for Heterogeneous CPU/GPU Programming (George Fekete)



Questions & Answers

NCCS User Services:

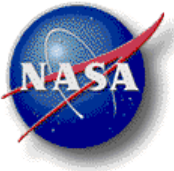
support@nccs.nasa.gov

301-286-9120

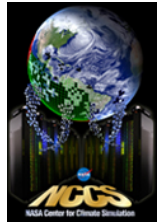
<https://www.nccs.nasa.gov>

NCCS User News Twitter feed at

http://twitter.com/NASA_NCCS



Contact Information



NCCS User Services:

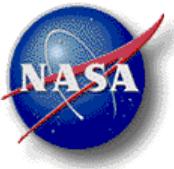
support@nccs.nasa.gov

301-286-9120

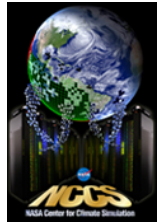
<https://www.nccs.nasa.gov>

NCCS User News Twitter feed at

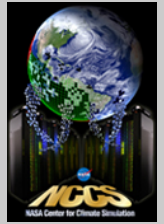
http://twitter.com/NASA_NCCS



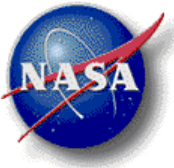
Agenda – March 20, 2012



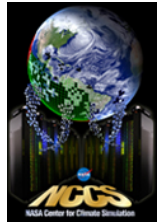
- Welcome & Introduction (Phil Webster or Lynn Parnell)
- Discover Update (Dan Duffy)
- Dirac DMF Growth Rate – Issue & Remediation (Tom Schardt)
- NCCS Operations & User Services (Ellen Salmon)
- Question & Answer
- Breakout Sessions
 - Overview of Debugging on Discover (Doris Pan)
 - OpenCL Framework for Heterogeneous CPU/GPU Programming (George Fekete)



Supporting Slides



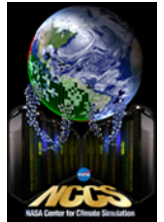
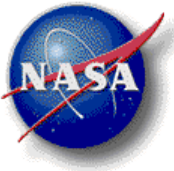
SMD Allocation Request URL (e-Books)



SMD Allocation Applications Deadline: 20 March (today!)

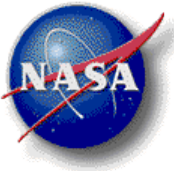
<http://hec.nasa.gov/request/announcements.html#spring2012>

<https://hec.reisys.com/hec/computing/index.do>

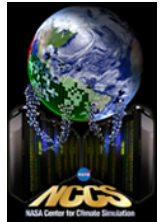


Supporting Slides – HPC

(Dan Duffy)

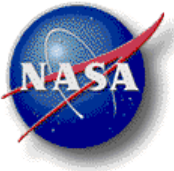


FY11 NCCS Upgrade of SCU3/SCU4

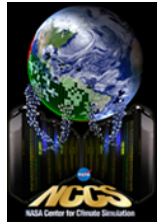


- Swap out of old IBM nodes installed in late 2008 with new IBM nodes
- Reused the infrastructure of the racks and switches
- 100% in production by October 2011
- Increase of ~2.6M SBUs per year

Feature	Old SCU3/SCU4	Upgrade SCU3+/SCU4+
Nodes	516	516
Processor Type	Intel Harpertown	Intel Westmere
Cores/Node	8	12
SBU Rate/Node	0.38	0.95
SBU Rate/Year	1,717,661	4,294,152

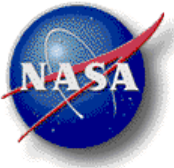


FY11 NCCS Upgrade of SCU1+/SCU2+

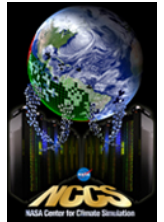


- Swap out of old Dell nodes installed in late 2006 and early 2007 with new IBM nodes
- Reused the infrastructure of the racks and switches
- 100% in production during month of January 2012
- Increase of ~3.4M SBUs per year

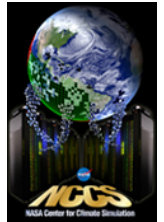
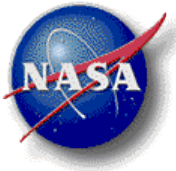
Feature	Old SCU1/SCU2	Upgrade SCU1+/SCU2+
Nodes	516	516
Processor Type	Intel Woodcrest	Intel Westmere
Cores/Node	4	12
SBU Rate/Node	0.2	0.95
SBU Rate/Year	904,032	4,294,152



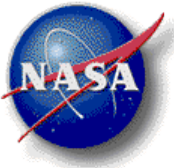
Dali-GPU Analysis Nodes Available



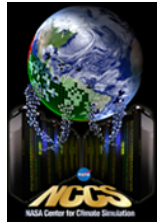
- 9 Dali-GPU Analysis nodes each configured with
 - Dual socket, hex-core Intel Westmere 2.8 GHz processors
 - 192 GB of RAM
 - QDR Infiniband
 - All GPFS file systems
 - Two Nvidia M2070 GPUs
 - 10 GbE Network Interface Cards
- Accessible by direct login
 - ssh dali.nccs.nasa.gov – might land on a node with a GPU (DNS round robin)
 - ssh dali-gpu.nccs.nasa.gov – gaurenteed to land on a node with a GPU
 - <http://www.nccs.nasa.gov/primer/getstarted.html#login>



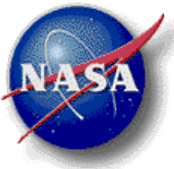
Supporting Slides – DMF



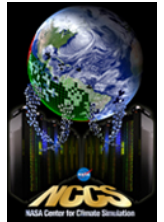
Costs



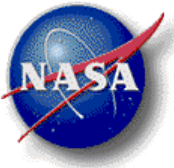
- DMF is licensed based on the amount of data managed.
 - A 5 PB Incremental License and three years of maintenance is \$62,500 (list)
- Addition tape media
 - 1000 Tapes (5 PBs) is \$300,000 (list)
- Costs of 1 year at 2 PBs per month growth is \$1,812,500 (list)
- Costs of moving off of old media is \$500,000 (list)



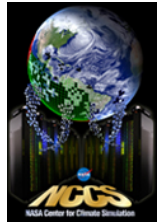
DMF Migration Process



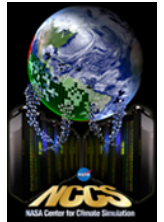
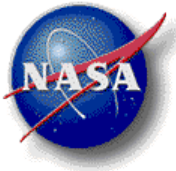
- DMF periodically scans each cache file system for new files – files appended to or overwritten are considered new
- A list of files to write to tape is created and run against the migration rules for the particular file system
- The number of tape copies written is determined by the migration rules



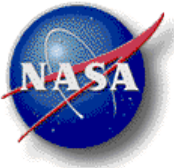
DMF Migration (cont.)



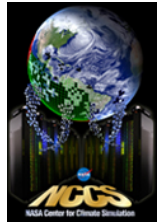
- The migration rules can be tailored by a file's age, group, owner, or "sitetag"
- The migration rules are only used when the file is initially written to tape
- Changing the number of tape copies of an existing DMF managed file will require updating the DMF database



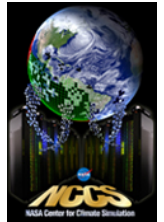
Supporting Slides – HPC Operations



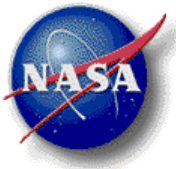
Chilled Water Outage Detail – March 13, 2012



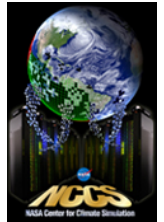
- NCCS system administrators brought down HPC systems in response to rapid, significant temperature increases during February 28 chilled water maintenance
 - Cluster unavailable to users most of that day
 - Restoring full access to archive system incrementally since event
 - Multiple heat-related failures in DDN storage devices
 - Ongoing failures anticipated
 - One large Discover filesystem, two archive filesystems still being recovered
 - Supported GMAO monthly forecast by temporarily mounting failed filesystem in read-only mode, GMAO copied necessary data to alternate location
 - Received numerous replacement disk drives from vendors. Installing and rebuilding disk storage tiers in cautious fashion to minimize data loss
 - Received 34 replacement disk drives for Discover and 14 for archive system [as of March 13, 2012]
 - Replaced 8 drives on Discover and 21 on archive system (had 7 spares for use on archive system) [as of March 13, 2012]



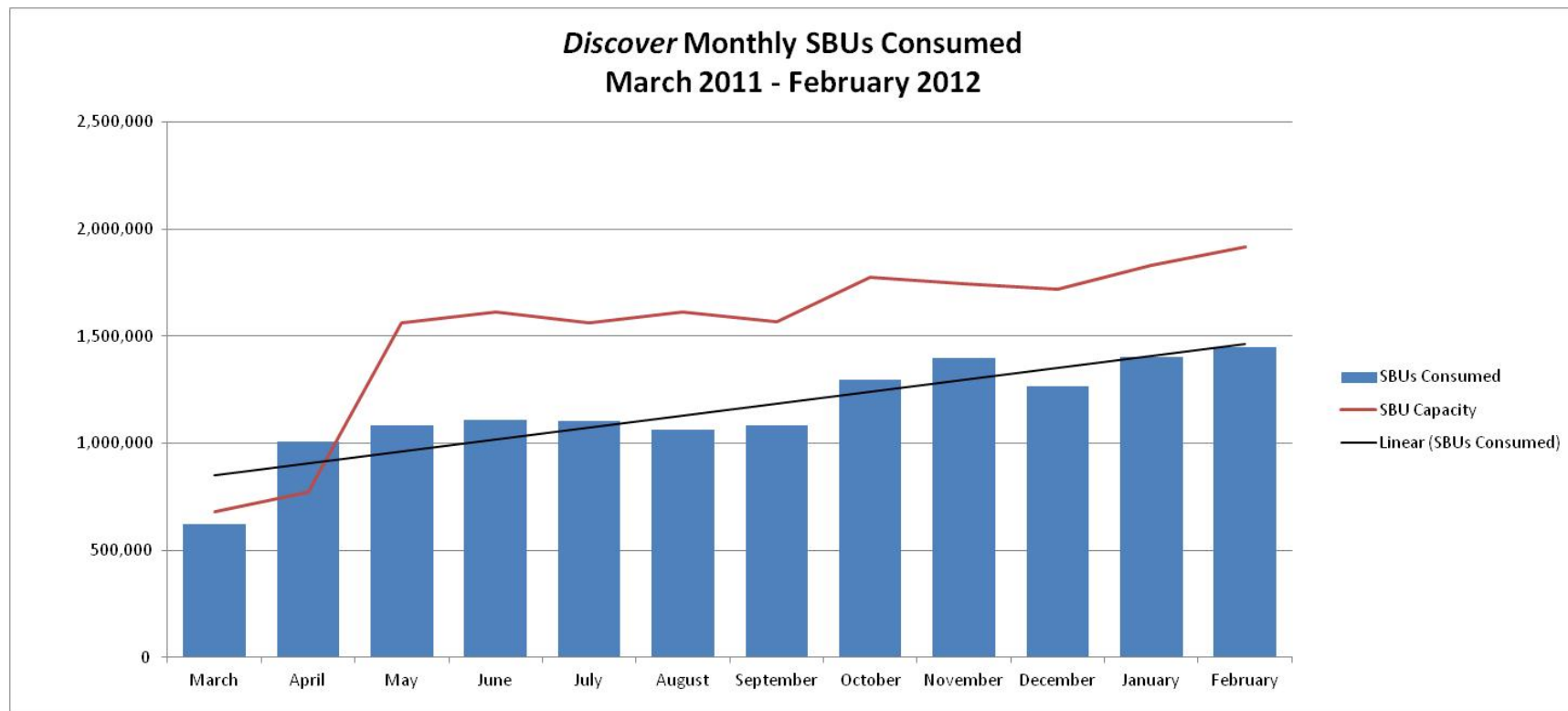
Supporting Slides – NCCS METRICS / Utilization

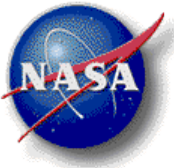


Discover Utilization



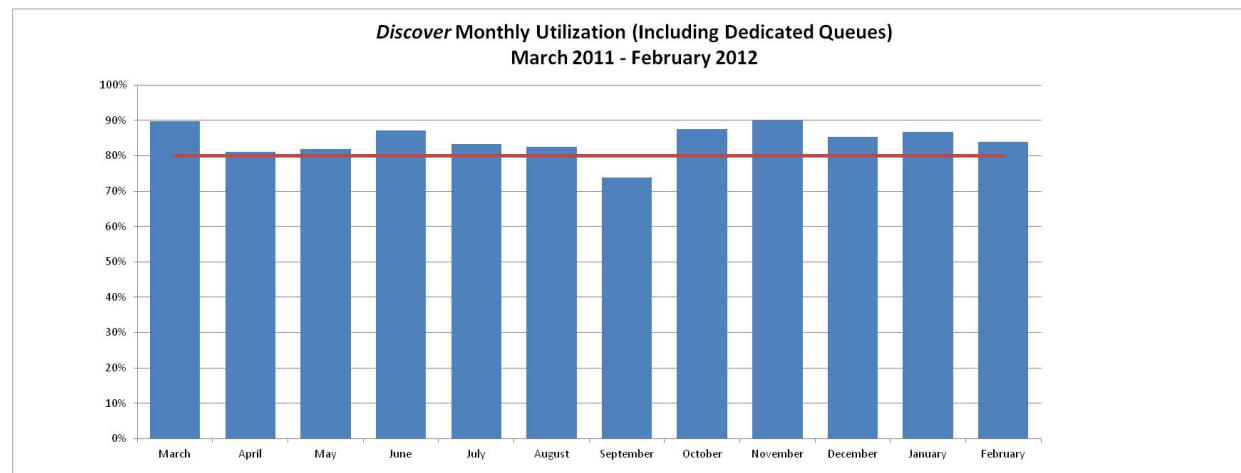
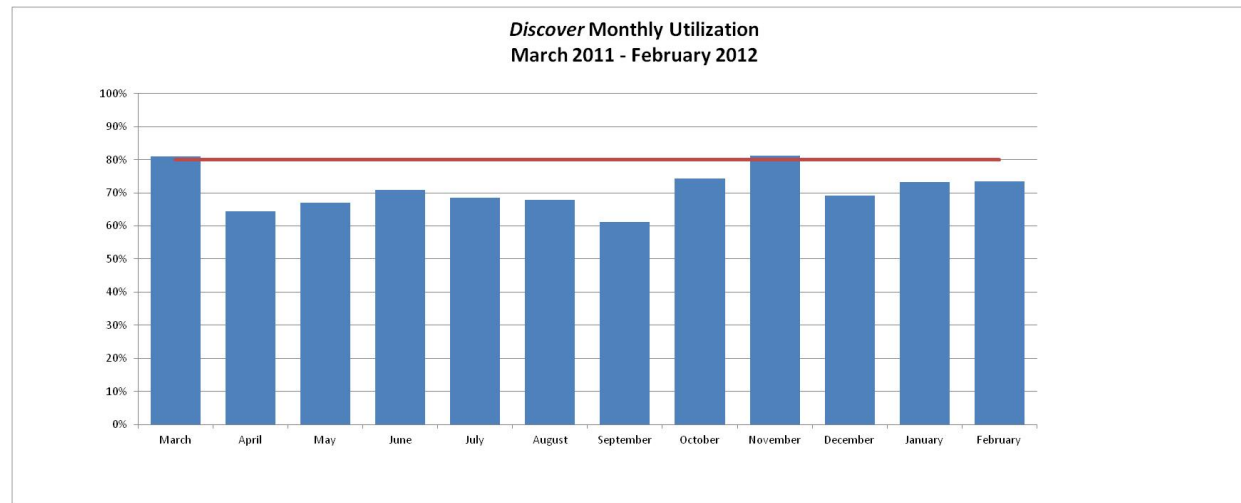
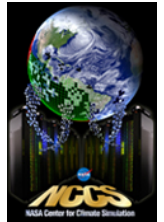
- Generally increasing SBU consumption
- 1.4 million SBUs consumed in February (about the same as January)
- 1.9 million SBUs theoretical maximum for February





Discover Utilization Percentage

- 80% target line
- Top chart based on actual SBU consumption
- Bottom chart treats dedicated queues as fully utilized



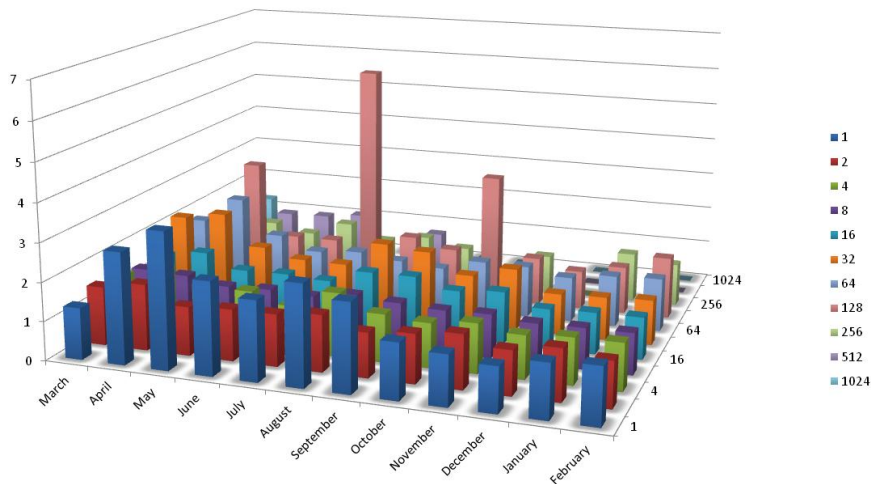


Discover Utilization and Wait

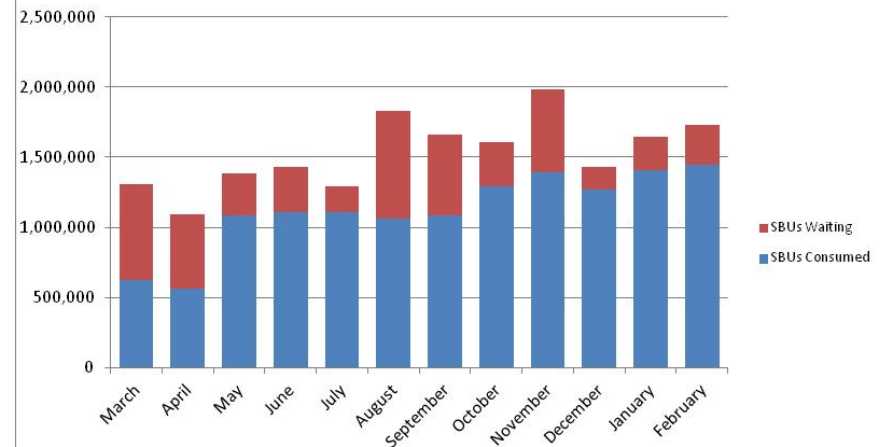


- Still observing latent demand (see red areas in stacked bar charts) – 283 thousand SBU equivalents spent awaiting job execution in February (242 thousand in January)
- Aggregate February expansion factor: 1.2 (about the same as January)
- February general_small queue expansion factor: 1.5 (large number of concurrent jobs submitted)
- February 128-node jobs' expansion factor: 1.62 (3x jobs submitted in this category compared to January)
- Earlier spikes generally due to user job submission patterns

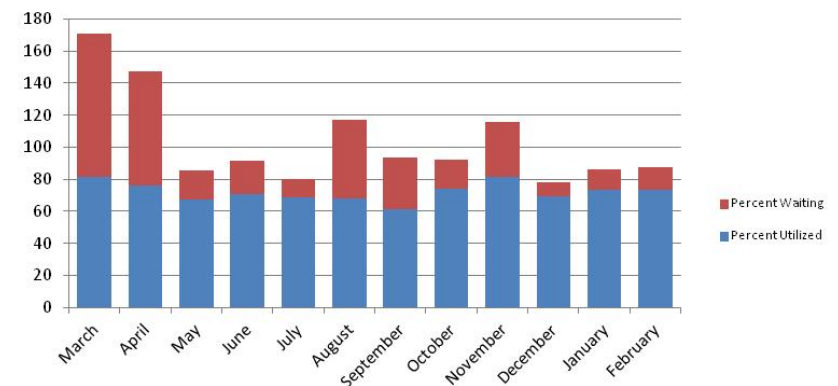
Discover Expansion Factors by Requested Nodes
March 2011 - February 2012

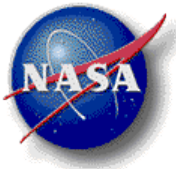


Discover SBUs Utilized and Awaiting Execution
March 2011 - February 2012

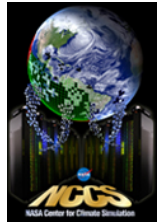


Discover Utilization and Wait Percentages
March 2011 - February 2012



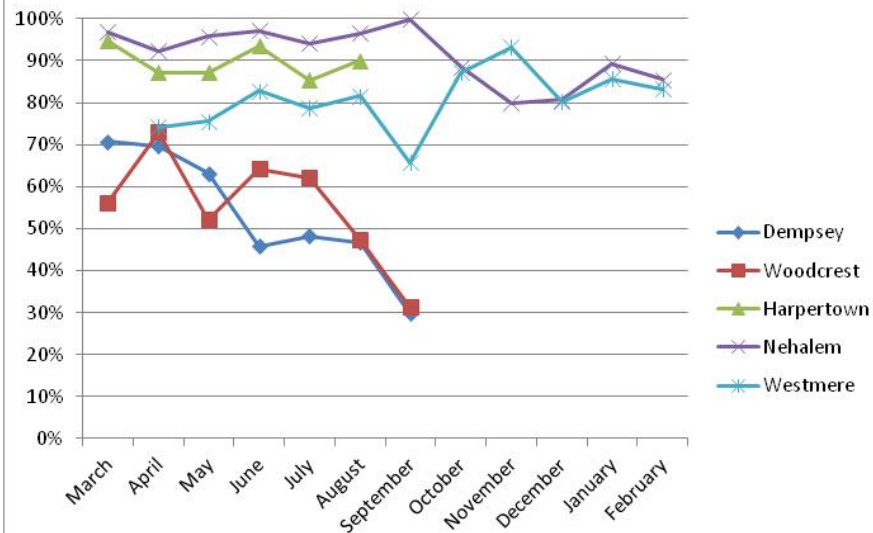


Discover Node Utilization



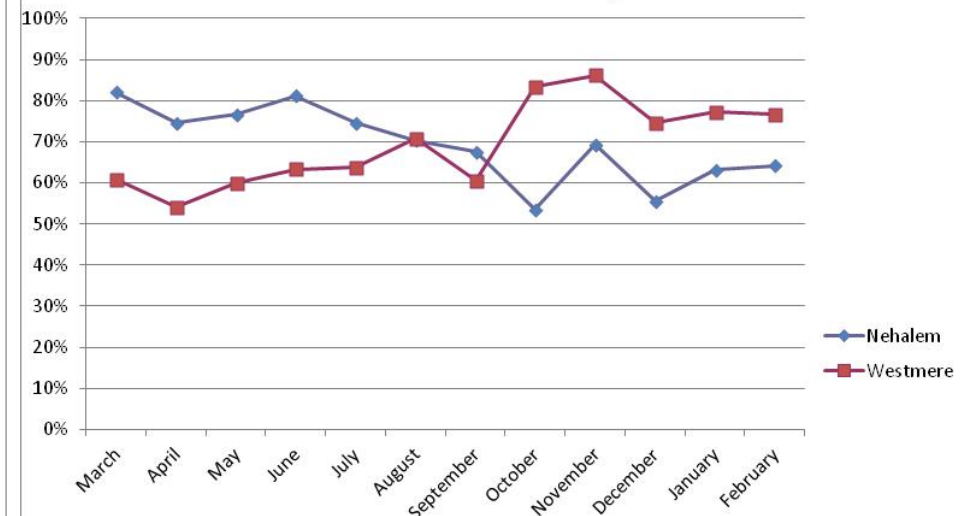
- Westmere and Nehalem both over 80% utilized (when factoring in dedicated nodes)
- Processor utilization lower when examining actual SBU consumption

Discover Processor Utilization by Architecture
March 2011 - February 2012



Assumes full utilization of dedicated queues

Discover Processor Utilization by Architecture
March 2011 - February 2012



Based on SBUs consumed

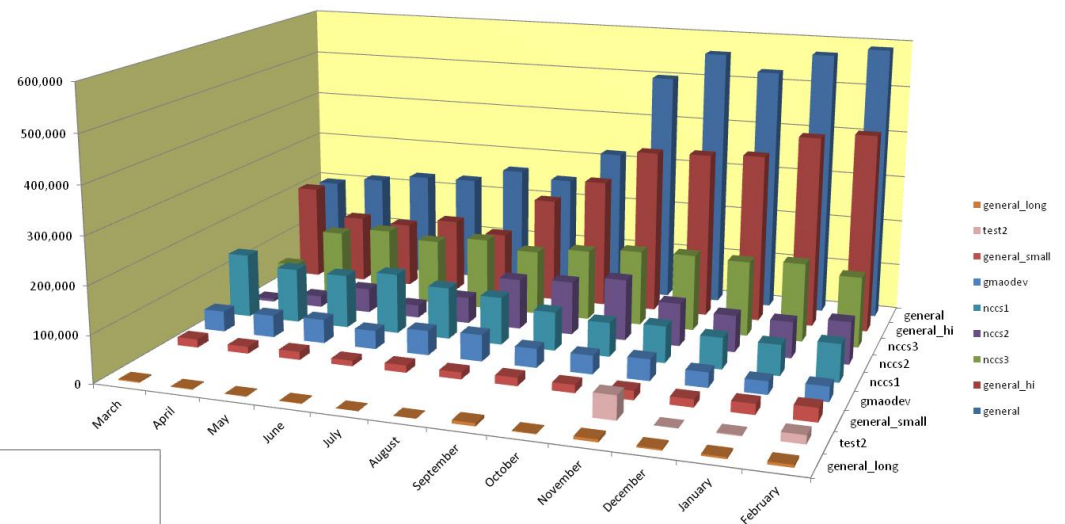


Discover SBU Distribution by Job Size and Queue

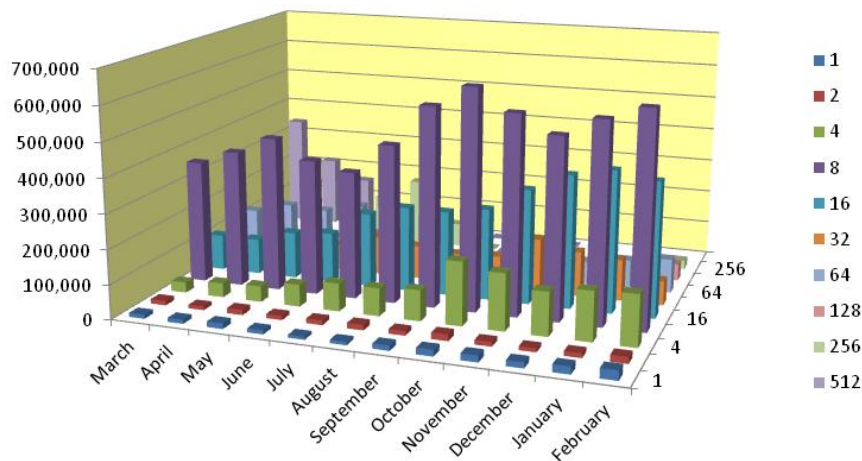


- 236 thousand jobs run in February (216 thousand in January)
- Consistent consumption pattern
- Most work performed in general, general_high queues
- Next greatest consumption observed in dedicated queues: nccs3, nccs2, nccs1

Discover SBUs by Top Queues
March 2011 - February 2012



Discover SBUs by Job Size (Nodes)
March 2011 - February 2012



- Consistent consumption pattern
- 8-node category (i.e., 8-15 nodes) used greatest number of SBUs (616 thousand in February, 577 thousand in January) by job size
- 16-node category (i.e., 16-31 nodes) next (390 thousand in February, 411 thousand in January)



Discover Job Distribution by Job Size and Queue

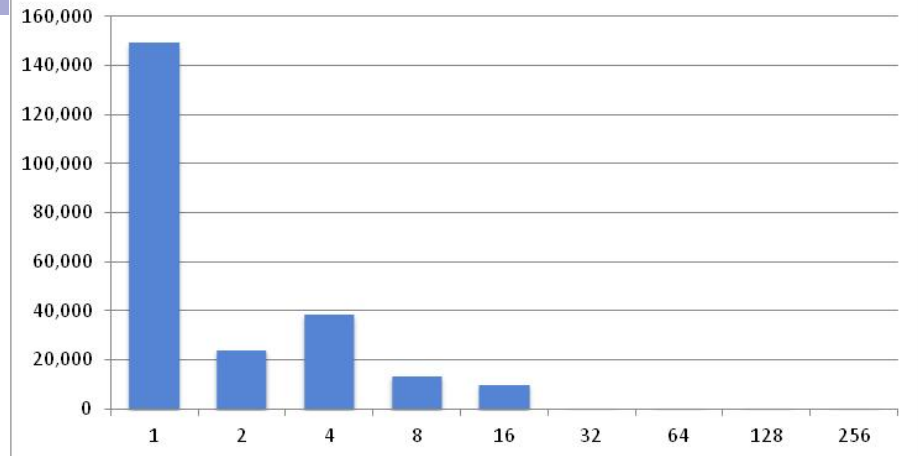


- 1-node category had greatest number of jobs by job size (typical)
- datamove queue had greatest number of jobs by queue (typical)

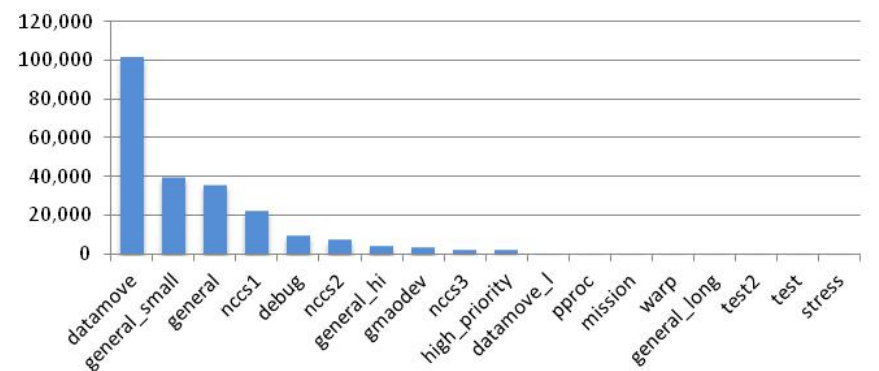
Nodes	Number of Jobs
1	149,204
2	23,895
4	38,498
8	13,504
16	9,818
32	419
64	521
128	97
256	23

Queue	Number of Jobs
datamove	101,536
general_small	39,557
general	35,350
nccs1	21,774
debug	9,358
nccs2	6,990
general_hi	3,406
gmaodev	2,745
nccs3	1,591
high_priority	1,362
datamove_l	721
pproc	529
mission	325
warp	107
general_long	42

**Discover Jobs By Job Size (Nodes)
February 2012**



**Discover Jobs by Queue
February 2012**





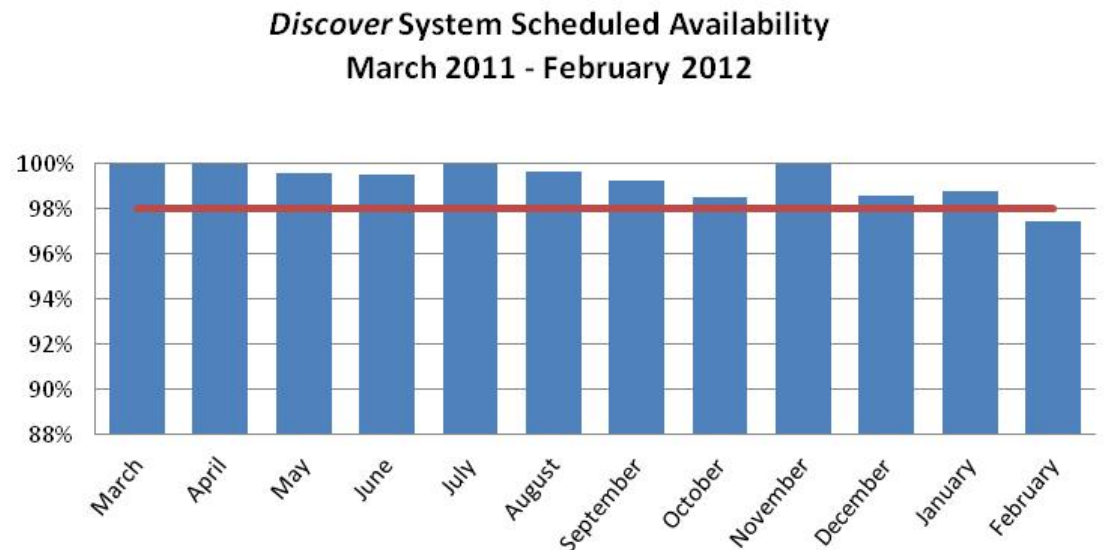
Discover System Availability



- April 2011 – System shut down for weekend due to potential furlough
- August 2011 – Two-week SCU3/SCU4 upgrade

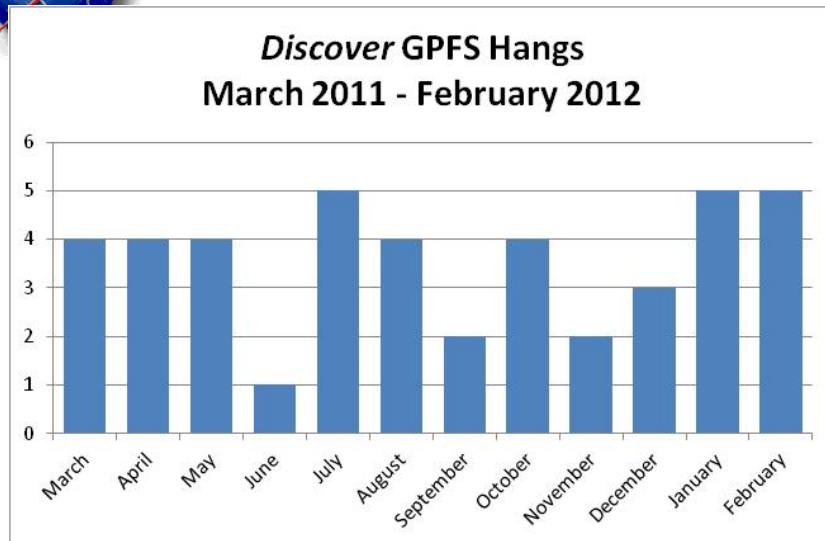
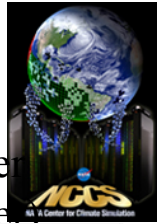


- Overall February system availability: 97.4% (compared to 98.8% in January)
- Experienced February outage due to high temperature resulting from chilled water outage



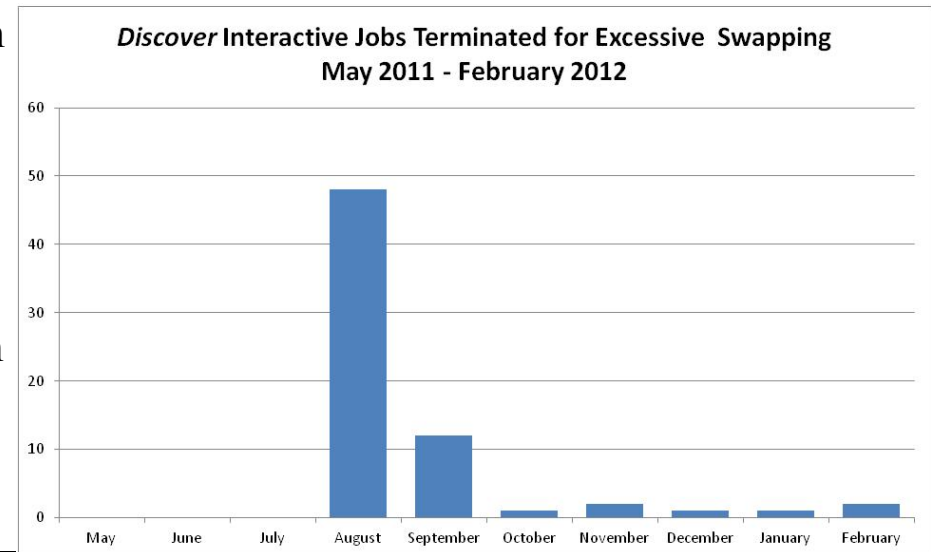
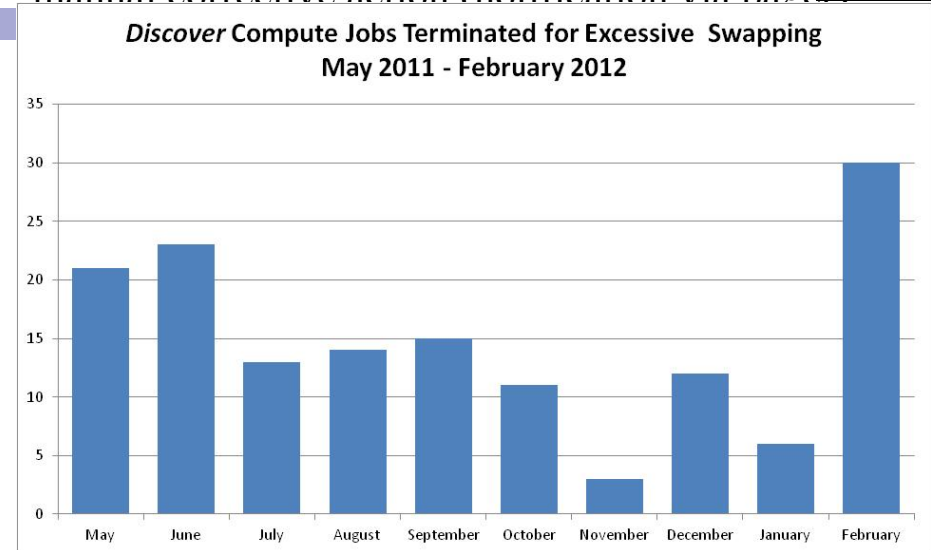


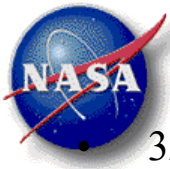
Discover Out of Memory Summary



- Compute node swap monitor script caused 38 compute node reboots in February (compared to 6 in January), terminating 30 different jobs (6 in January) and affecting 14 distinct users (6 in January)
- Decreased threshold in February from 80% to 70% swap
- Interactive node (Dali and login) swap monitor script terminated 2 interactive users' processes in February (1 in January)

- Five GPFS hangs of sufficient impact to trigger manual corrective action (notification via page)

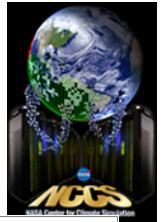




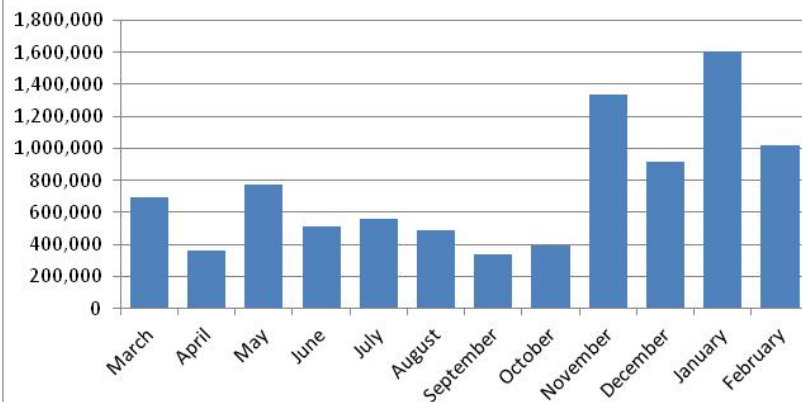
Mass Storage Utilization

32.0.PB (Base 10) in library (compared to 29.8 PB in January)

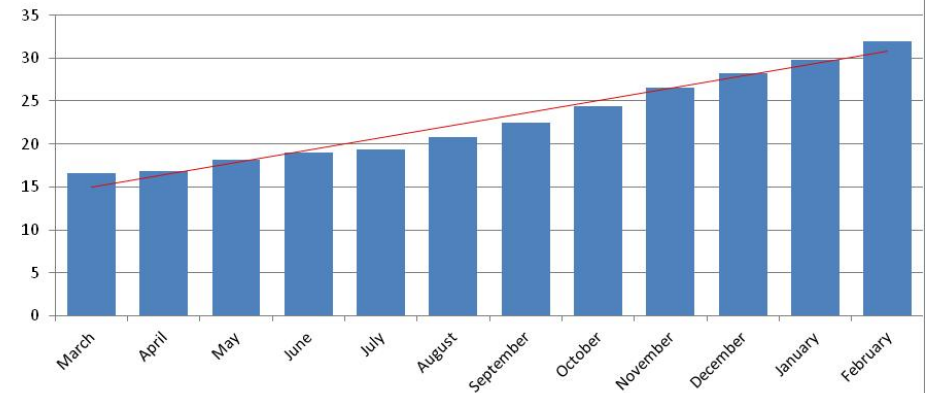
- When looking at library holdings since mid-2006, recent rate of increase appears to be accelerating much more rapidly
- 1,016,808 DMF recalls (user requests) in February (compared to 1,606,091 in January)



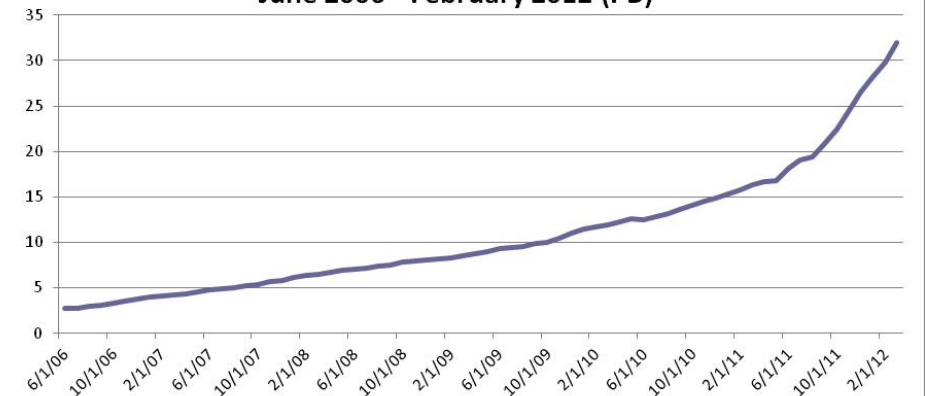
Archive Recalls
March 2011 - February 2012



Mass Storage Data Stored
March 2011 - February 2012 (PB)

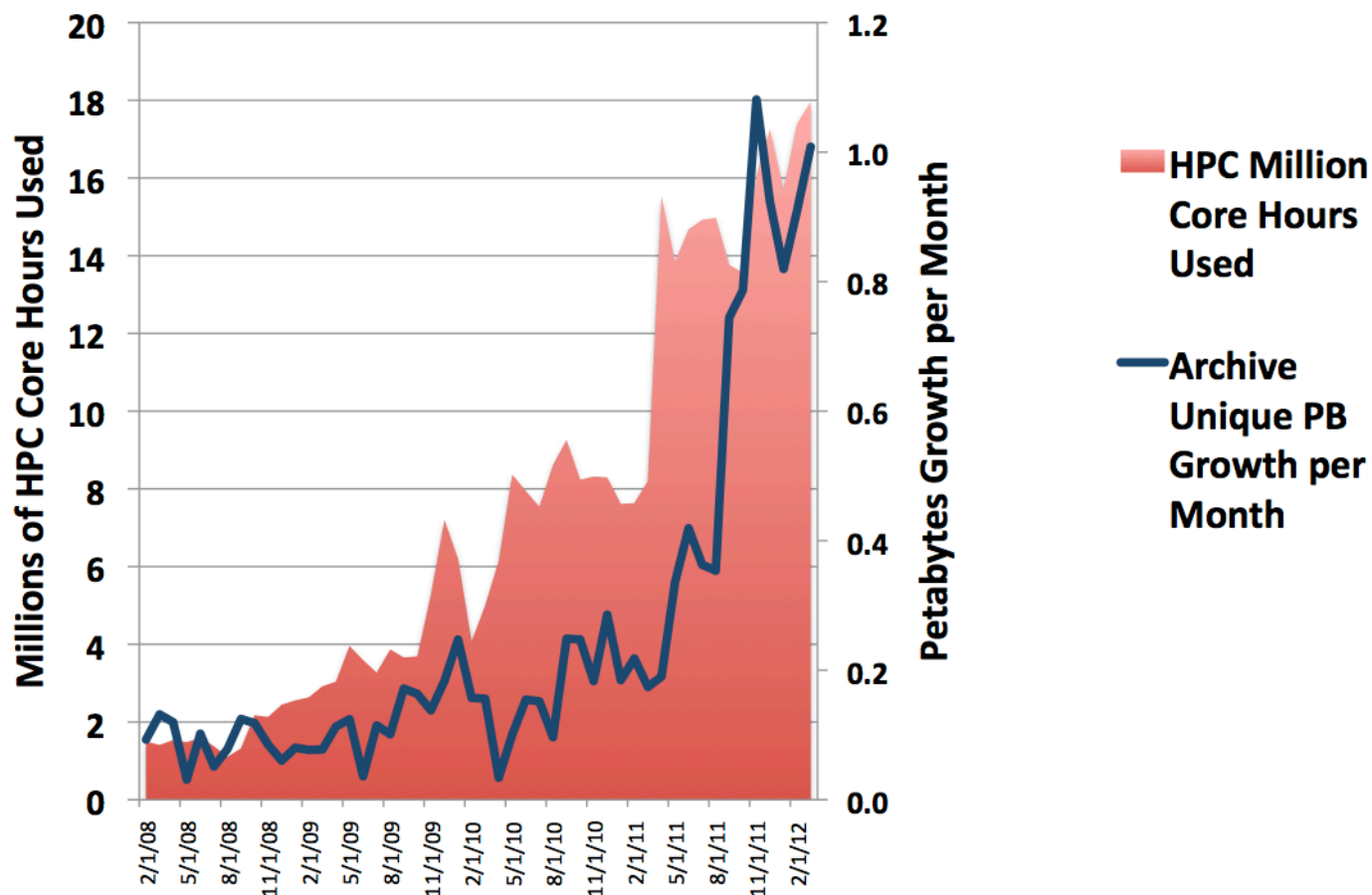


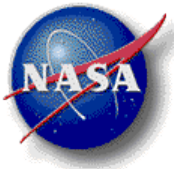
Mass Storage Data Stored
June 2006 - February 2012 (PB)





NCCS Archive Growth per Month and CPU Hours Utilized, February 2008 - February 2012





NCCS Monthly Archive Growth (Single Copy) and System Billable Units (SBUs), July 2010 - February 2012

